

# A Appendix

## A.1 Interaction Design

*SemiConLens* provides a set of interactions to support the design tasks. The interactions are designed to be consistent across different views, enabling users to seamlessly transition between views and maintain a coherent mental model for the sequential exploration. The interactions are as follows:

**Lasso.** Users can draw a lasso selection in the Discovery View to select multiple compounds for cluster and reference filtering.

**Brushing.** Users can use brush in the Attribute Distribution of the Filter View to define value ranges for range filtering.

**Clicking.** Users can click on a glyph in the Discovery View to access detailed compound information in the Comparison View. Additionally, clicking on a filter state in the Exploration History navigates users to trace back and forth through the filtering process.

**Highlighting.** Users can hover over a glyph in the Comparison View to highlight the corresponding compound in the Discovery View and Filter View. Additionally, hovering over a filter state in the Exploration History highlights the corresponding attribute distribution in the Attribute Distribution.

**Tooltip.** *SemiConLens* provides tooltips in the Discovery View and Comparison View to display detailed attribute information upon hovering over a specific glyph or heatmap cell. The tooltip includes the exact attribute value and uncertainty.

## A.2 Data sources and preprocessing

Several major databases provide material attributes. The C2DB<sup>1</sup> reports computational properties of 16,905 2D materials, the Materials Project<sup>2</sup> covers 169,385 compounds, 2DMatPedia<sup>3</sup> compiles 6,351 entries from theory and experiment, and eTran2D<sup>4</sup> focuses on transport data for 91 compounds.

We used C2DB as the primary source due to its size and consistent DFT framework, merging it with mobility values from eTran2D. The combined data were stored in MySQL and filtered (0.2- 4.0 eV band gap, dynamic stability, real phonons, non-magnetic). This yielded 1,339 candidates, with mobility values for 50 compounds and GW-calculated gap, CBM, VBM for 225.

---

<sup>1</sup><https://2dhub.org/c2db/c2db.html>

<sup>2</sup><https://next-gen.materialsproject.org/>

<sup>3</sup><http://www.2dmatpedia.org/>

<sup>4</sup><https://sites.utexas.edu/yuanyue-liu/etran2d/>

### A.3 The System Workflow

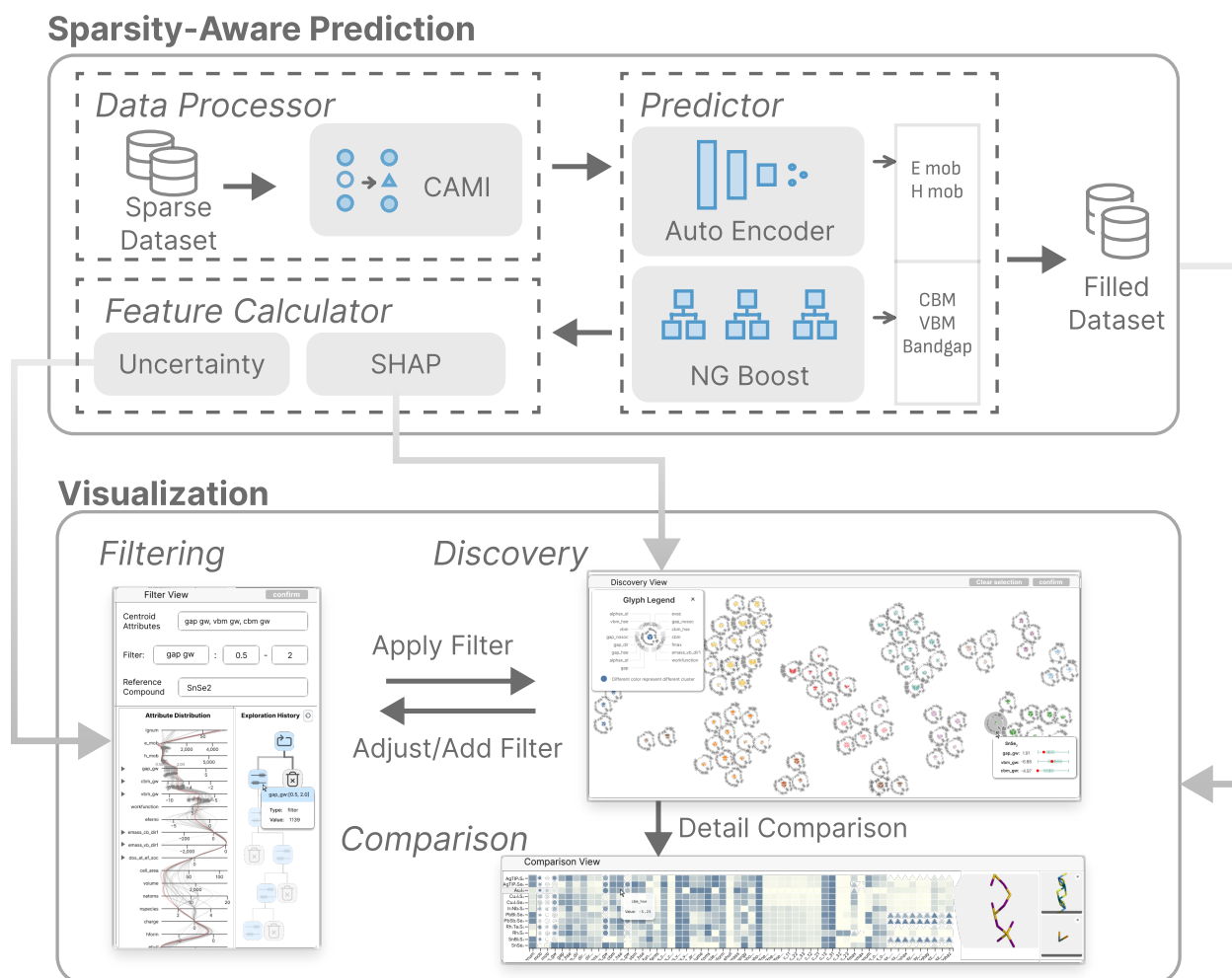


Fig. 1. The workflow consists of two main components: (1) a sparsity-aware prediction pipeline that processes sparse datasets through CAMI data processing, followed by prediction using Auto Encoder and NG Boost models, with uncertainty and SHAP-based feature calculation for interpretability; (2) a visual analysis system, which enables interactive *filtering*, *exploration* and *comparison* of the filled dataset.

## A.4 AE Finetuning

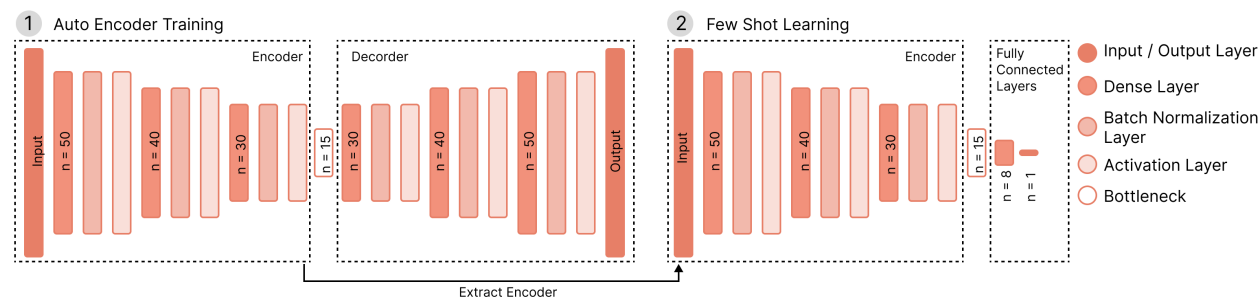


Fig. 2. Configuration of the Auto Encoder and few shot learning model.

As discussed in Quantitative Evaluation in section 7.1, Fig. 2 demonstrates the architecture of the AE in the training and few shot learning. In step Fig. 2 1, an AE was trained on all 1,339 compounds. Then the encoder and the bottleneck (latent space) was extracted and the decoder was replaced with a small fully connected network as depicted in Fig. 2 2. For AE training the total number of epochs trained were 1000 with an early stopping mechanism of patience value of 10 on validation loss. Similar for AE few shot learning the patience value was 10 and number of epochs was 200.

## A.5 Questions

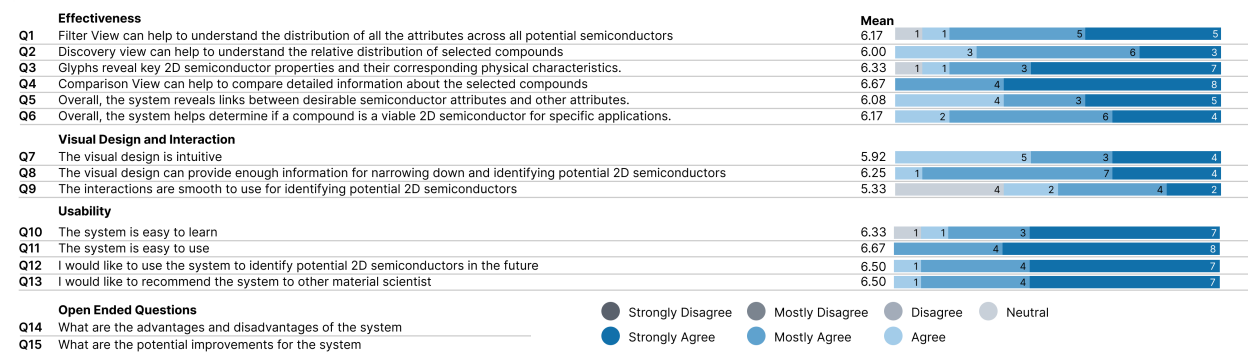


Fig. 3. The expert interview questionnaire results. Q1 to Q13 are close-ended questions which were answered on a 7-point Likert scale. Q14 and Q15 are open-ended questions.

Fig.3 reports the unbridged post study questions we used for the study with the summary of responses from the experts.

## A.6 Pseudo code for the Proposed Correlation Aware Multivariate Imputation

---

### Algorithm 1 Correlation Aware Multivariate Imputation (CAMI)

---

**Input:** Dataset  $X \in \mathbb{R}^{m \times n}$ , thresholds  $A$  and  $B$  (in percentages), and parameters  $C$  (number of correlated columns to consider) and  $D$  (number of similar rows to consider).

- 1:  $X \leftarrow X_{:,C}$  where  $C = \{j \in \{1, \dots, n\} \mid \frac{| \{i \mid X_{ij} = \text{null}\} |}{m} \leq A\%\}$
- 2:  $X \leftarrow X_{\mathcal{R},:}$  where  $\mathcal{R} = \{i \in \{1, \dots, m\} \mid \frac{| \{j \mid X_{ij} = \text{null}\} |}{|C|} \leq B\%\}$
- 3:  $c_m \leftarrow \text{Corr}(X)$
- 4: **for** each column  $Y$  in  $X$  **do**
- 5:     Let  $j$  be the index of  $Y$ .
- 6:      $K(Y) \leftarrow$  the set of top  $C$  indices from  $\{1, \dots, n\} \setminus \{j\}$  based on the correlation values in  $c_m$ .
- 7:     **for** each  $k \in K(Y)$  **do**
- 8:         Compute  $\mu_k = \frac{1}{|\mathcal{I}_k|} \sum_{i \in \mathcal{I}_k} X_{ik}$ , where  $\mathcal{I}_k = \{i \mid X_{ik} \neq \text{null}\}$ .
- 9:         Replace each missing value in column  $k$  with  $\mu_k$ .
- 10:     **end for**
- 11:     Let  $X' = X_{:,\{j\} \cup K(Y)}$ .
- 12:     Define  $\mathcal{L} = \{i \mid Y_i = \text{null}\}$  and  $\mathcal{M} = \{i \mid Y_i \neq \text{null}\}$ .
- 13:     **for** each row  $P \in \mathcal{L}$  **do**
- 14:         Compute the similarity between  $P$  and every row  $Q \in \mathcal{M}$  using features in  $K(Y)$ .
- 15:         Let  $\mathcal{O}_P$  be the top  $D$  rows in  $\mathcal{M}$  that are most similar to  $P$ .
- 16:         Impute the missing value in  $P$  by  $Y_P \leftarrow \frac{1}{D} \sum_{Q \in \mathcal{O}_P} Y_Q$ .
- 17:     **end for**
- 18: **end for**

**Output:**  $X$

---

Hyperparameters were iteratively tuned for the dataset used. Hyperparameter values were A=40%, B=40%, C=10 and D=10.

## A.7 An illustration of the proposed Correlation Aware Multivariate Imputation

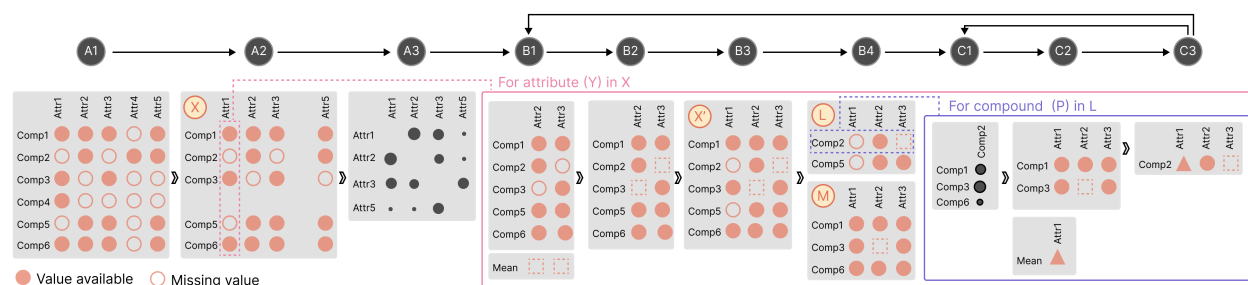


Fig. 4. An illustration of the proposed Correlation Aware Multivariate Imputation. Step **A1** : Load the dataset. Step **A2** : Select only columns and rows with less than A% and B% missing values respectively. Step **A3** : Calculate the correlation matrix. Step **B1** : For each attribute (Y) in the filtered dataset X, select top C most correlated columns using the correlation matrix. Then calculate the mean for each selected column. Step **B2** : Fill the missing values of each column with the corresponding calculated mean. Step **B3** : Create X', a subset of X, using the Y and its most correlated columns mentioned above. Step **B4** : split X' into L and M, where column Y is missing and not missing respectively. Step **C1** : For each compound (P) in L, compute similarity matrix with compounds in M. Step **C2** : Select the top D most similar compounds from the subset and calculate the mean for column Y. **C3** : Fill the missing value of P with the mean calculated above.



## A.9 Case 1 Exploring clusters - Identifying 2D Materials for Water Splitting

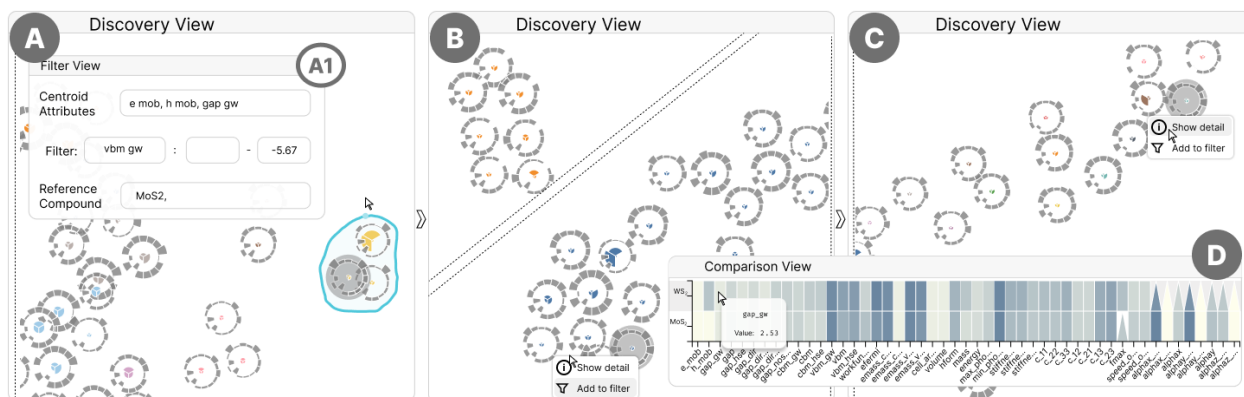


Fig. 6. With *SemiConLens* an expert has identified WS<sub>2</sub> as a similar material to MoS<sub>2</sub>, which is suitable for hydrogen evolution reaction in water splitting. First, the expert applied 2 filters, CBM GW > -4.44eV and VBM GW < -5.67 eV using the Filter View. Then they selected MoS<sub>2</sub> as the reference compound as shown in (A1). From the resulting Discovery View, they identified the cluster with MoS<sub>2</sub> and selected it by lassoing the cluster as depicted in (A). They identified the MoS<sub>2</sub> from the resulting Discovery View (B) and add it as the filter for next step. Then they identified WS<sub>2</sub> as the nearest neighbor for MoS<sub>2</sub> from the Discovery View (C) and added both to comparison view for further exploration as shown in (D). Some compounds are excluded from the Discovery View, as indicated by the dotted lines.

## A.10 Space Optimization for Cluster-Aware Overlap Removal

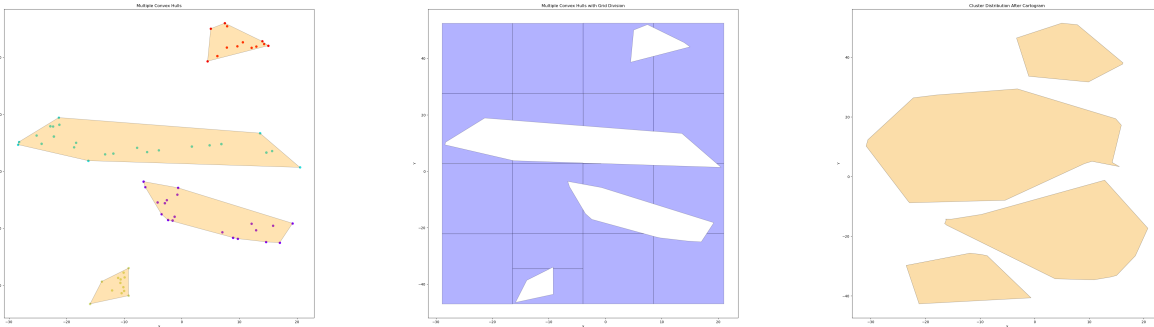


Fig. 7. Examples of intermediate representations in the space-optimization pipeline: (a) cluster polygons constructed from 2D point embeddings; (b) dummy regions generated to fill the remaining space; and (c) the final deformed layout after cartogram transformation.

**Cluster Region Construction.** The first step transforms discrete point clusters into continuous spatial regions while preserving cluster relationships. Specifically, each cluster is represented as a convex hull polygon  $P_i$ , computed using the Quickhull algorithm on the 2D embedding space. For clusters with fewer than three points, we instead generate a circular buffer of fixed radius around the centroid to avoid degenerate polygons.

**Dummy Region Generation.** Because clustering algorithms often produce non-adjacent cluster polygons, dummy polygons are introduced to fill the empty spaces for subsequent cartogram processing. These dummy regions are generated through a kd-tree inspired recursive partitioning procedure. We first compute a global bounding box covering all cluster polygons and initialize it as a  $4 \times 4$  grid.

**Recursive Subdivision and Clipping.** For each grid cell that overlaps the remaining empty space, if the cell fully contains a cluster polygon, we recursively split it along its longer axis until its area falls below 0.1 in normalized units. The resulting regions are then clipped to the empty space, balancing computational efficiency and geometric fidelity.

**Weight Assignment and Cartogram Transformation.** With this continuous layout in place, we then employ the diffusion-based cartogram technique to optimize polygon areas. Each cluster polygon is assigned a weight proportional to its number of member compounds, while dummy polygons are assigned a small constant weight of 0.1 so that they act only as spacers. This ensures sufficient room for expansion while approximately preserving relative topological relationships.

**Pseudo-code.**

---

**Algorithm 2** Cluster-Aware Space Optimization

---

**Require:** Points  $P = \{(x_i, y_i, c_i)\}$  in 2D embedding space

**Ensure:** Deformed polygons  $G$

- 1: Group points by cluster label
- 2: **for** each cluster  $c_i$  **do**
- 3:     **if**  $|P_i| < 3$  **then**
- 4:         Construct a circular buffer around the cluster centroid
- 5:     **else**
- 6:         Compute the convex hull polygon  $P_i$  using Quickhull
- 7:     **end if**
- 8: **end for**
- 9: Compute the global bounding box  $B$  covering all cluster polygons
- 10: Compute the remaining empty space  $R \leftarrow B \setminus \bigcup_i P_i$
- 11: Initialize a  $4 \times 4$  grid over  $B$
- 12: **for** each grid cell  $g$  **do**
- 13:     **if**  $g$  intersects  $R$  **then**
- 14:         **if**  $g$  fully contains a cluster polygon **then**
- 15:             Recursively subdivide  $g$  along its longer axis
- 16:             Stop subdivision when cell area  $< 0.1$
- 17:         **end if**
- 18:         Clip the resulting cell(s) to  $R$  and add them as dummy polygons
- 19:     **end if**
- 20: **end for**
- 21: Assign cluster weights according to cluster sizes
- 22: Assign dummy polygons a small constant weight of 0.1
- 23:  $G_0 \leftarrow$  cluster polygons  $\cup$  dummy polygons
- 24:  $G \leftarrow$  Cartogram( $G_0, w$ )
- 25: **return**  $G$

---